

**Supporting Information**

**for**

**Analysis of Protein Dynamics Simulations**

**by a Stochastic Point Process Approach**

Bertil Halle and Filip Persson

Biophysical Chemistry, Lund University, POB 124, SE-22100 Lund, Sweden

## S1. CODE FOR RT HISTOGRAM

The following Matlab code computes the RT histogram  $\mathbf{F}_R$ , given the occupancy vector  $\mathbf{A}_0$  (including any vacancies) for a singly occupied site.

```
A = A0(A0~=0);  
  
B = any(diff(A,1,1)~=0,2);  
  
B = [true; B(:); true];  
  
VR = diff(find(B));  
  
if A0(length(A0))~=0, VR(length(VR)) = []; end  
  
if A0(1)~=0, VR(1) = []; end  
  
FR = zeros(max(VR),1);  
  
VRuni = unique(VR);  
  
FR(VRuni) = histc(VR,VRuni);
```

## S2. CODE FOR RT AND ST STATISTICS

The following Matlab code computes the quantities  $Q_R(n)$ ,  $Q_S(n)$ ,  $\tau_R$  and  $\tau_S$ , given the RT vector  $\mathbf{V}_R$ , the RT histogram  $\mathbf{F}_R$  and the sampling resolution  $\Delta\tau$ . In addition, the uncorrelated statistical uncertainties in these quantities are computed.

```

Nvec = [0; (1:length(FR))'];

tau = dtau*Nvec;

NR = sum(FR);

NF = sum(VR);

VRm = NF/NR;

VR2m = mean(VR.^2);

VR3m = mean(VR.^3);

VR4m = mean(VR.^4);

tauR = dtau*VRm;

tauS = (dtau/2)*VR2m/VRm;

FRcumsum1 = flipud(cumsum(flipud(FR)));

QR = [1; (FRcumsum1 - FR)/NR];

FRcumsum2 = flipud(cumsum(flipud(FRcumsum1)));

QS = [1; (FRcumsum2 - FRcumsum1)/NF];

sigtauR = dtau*sqrt((VR2m - VRm^2)/(NR-1));

sigtauS = dtau*sqrt((VRm^2*VR4m + VR2m^3 - 2*VRm*VR2m*VR3m)/NR)/(2*VRm^2);

sigQR = sqrt(QR.*(1-QR)/(NR-1));

NQR = (2*Nvec+1).*QR;

sigQS = (dtau/tauR)*sqrt((cumsum(NQR) - NQR - (cumsum(QR) - QR).^2)/NR);

```

### S3. MULTIPLE OCCUPANCY

Here we indicate how the algorithm for constructing the RT vector can be generalized to multiply occupied hydration sites.

For a site that contains  $\nu$  water molecules at all times, we replace the occupation vector by an occupation matrix  $\mathbf{A}$  of size  $N_F \times \nu$ . To analyze the overall RT and ST statistics for the entire site, we need not keep track of the individual subsites or of any water interchange among them. Exchange events are identified by taking  $\nu \times \nu$  differences of  $w$  indices between successive frames. The site-based RTs obtained in this way are then multiplied by  $\nu$  to obtain molecule-based RTs averaged over the  $\nu$  subsites.

If the occupancy fluctuates in time, the  $\mathbf{A}_0$  matrix is dimensioned according to the maximum occupancy  $\nu_{\max}$  and a negative integer is entered instead of the molecule index  $w$  whenever one or more subsites are vacant (for example,  $-1$  and  $-2$  if two subsites are vacant in the same frame). The occupancy in a given frame  $k$  then equals the number of positive elements in the corresponding row of  $\mathbf{A}_0$ . Rather than removing all vacancies, as we did for the singly occupied site, we can treat the negative  $w$  indices as additional water molecules. The first-frame index must now be accompanied by a label that specifies whether it refers to a vacancy or not, so that the corresponding vacancy RTs can be discarded.

## S4. BINNING ERROR

Here we provide more details about the analysis of the binning error in the RT histogram  $F_R(n)$  and in the RT and ST statistics.

Figure 1 shows a small part of the analyzed trajectory. On the upper continuous time line, exchange events are indicated by tick marks and the different resident water molecules are labeled by letters. On the lower discrete time line, the observation time points, that is, the frames saved for analysis, are indicated by dots separated by the time interval  $\Delta\tau$ , which is the resolution. For each frame, the identity of the observed resident water molecule is indicated.

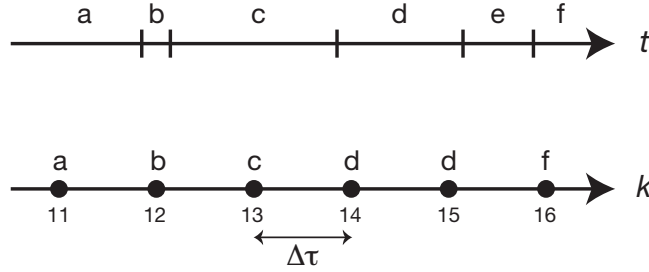


Figure S1: Top: continuous time line with 6 resident water molecules (labelled  $a$  to  $f$ ) and 5 exchange events (tick marks). Bottom: discrete time line indicating the water molecules observed at resolution  $\Delta\tau$ .

The distribution of RTs on the upper time line, that is, the separation of adjacent tick marks, is described by the continuous probability density  $\psi_{\text{R}}(\tau)$ , such that  $\psi_{\text{R}}(\tau) d\tau$  is the fraction of all RTs that are within  $d\tau$  of  $\tau$  in length. The total number of RTs in the trajectory is denoted by  $N_{\text{R}}^0$ .

From the lower time line, we obtain the observed RT histogram  $F_{\text{R}}(n)$ , which gives the number of times that we observe in the trajectory the same resident water molecule in precisely  $n$  contiguous frames. The total number of such discretized RTs is denoted by  $N_{\text{R}}$ . Clearly,

$$N_{\text{R}} = \sum_{n=1}^{\infty} F_{\text{R}}(n) \quad (\text{S1})$$

In Fig. S1,  $N_{\text{R}}^0 = 6$  but  $N_{\text{R}} = 5$  since water molecule  $e$  escapes detection.

From the example in Fig. S1, it is seen that water molecules  $b$  and  $c$  are both assigned a discrete RT of  $n = 1$  although the true continuous RTs are only a fraction of  $\Delta\tau$  or nearly  $2\Delta\tau$ , respectively. In general, if the discrete RT is  $n$ , we know that the continuous RT must lie in the interval

$$(n-1)\Delta\tau < \tau \leq (n+1)\Delta\tau \quad (\text{S2})$$

But not all continuous RTs in this interval contribute to  $F_{\text{R}}(n)$ . The mapping from the continuous probability density  $\psi_{\text{R}}(\tau)$  to the histogram  $F_{\text{R}}(n)$  takes the form

$$F_{\text{R}}(n) = N_{\text{R}}^0 \int_{(n-1)\Delta\tau}^{(n+1)\Delta\tau} d\tau \psi_{\text{R}}(\tau) p_n(\tau) \quad , \quad n \geq 1 \quad (\text{S3})$$

where  $N_{\text{R}}^0 \psi_{\text{r}}(\tau) d\tau$  is the number of continuous RTs with length within  $d\tau$  of  $\tau$  and  $p_n(\tau)$  is the probability that a randomly placed time interval of length  $\tau$  covers exactly  $n$  frames. Clearly,

$$p_n(\tau) = \begin{cases} \frac{\tau - (n-1)\Delta\tau}{\Delta\tau} & , \quad \text{for} \quad (n-1)\Delta\tau < \tau \leq n\Delta\tau \\ \frac{(n+1)\Delta\tau - \tau}{\Delta\tau} & , \quad \text{for} \quad n\Delta\tau < \tau \leq (n+1)\Delta\tau \end{cases} \quad (\text{S4})$$

Thus,

$$F_{\text{R}}(n) = \frac{N_{\text{R}}^0}{\Delta\tau} \left\{ \int_{(n-1)\Delta\tau}^{n\Delta\tau} d\tau \psi_{\text{R}}(\tau) [\tau - (n-1)\Delta\tau] + \int_{n\Delta\tau}^{(n+1)\Delta\tau} d\tau \psi_{\text{R}}(\tau) [(n+1)\Delta\tau - \tau] \right\} \quad (\text{S5})$$

For  $\tau$  in the interval specified in eq S2, we can Taylor expand  $\psi_{\text{R}}(\tau)$  around the central frame at  $\tau = n\Delta\tau$  as

$$\psi_{\text{R}}(\tau) = \psi_{\text{R}}(n\Delta\tau) + (\tau - n\Delta\tau) \left. \frac{\partial \psi_{\text{R}}}{\partial \tau} \right|_{\tau=n\Delta\tau} + \frac{(\tau - n\Delta\tau)^2}{2} \left. \frac{\partial^2 \psi_{\text{R}}}{\partial \tau^2} \right|_{\tau=n\Delta\tau} + \dots \quad (\text{S6})$$

Inserting this expansion into eq S5 and carrying out the integrals, we find that the linear term in eq S6 does not contribute and

$$F_{\text{R}}(n) = N_{\text{R}}^0 \Delta\tau \psi_{\text{R}}(n\Delta\tau) \left[ 1 + \mathcal{O}\left(\frac{\Delta\tau}{\tau_{\text{R}}}\right)^2 \right] \quad , \quad n \geq 1 \quad (\text{S7})$$

where, in the quadratic term, we have assumed that  $\psi_{\text{R}}''(\tau) = \psi_{\text{R}}(\tau)/\tau_{\text{R}}^2$ , as for a Poisson process.

The discretized correlation functions  $Q_{\text{X}}(n)$  are most conveniently expressed, not in terms of the observed histogram  $F_{\text{R}}(n)$ , but in terms of an ‘ideal’ histogram  $F_{\text{R}}^0(n)$

defined, for  $n \geq 1$ , as

$$F_{\text{R}}^0(n) = N_{\text{R}}^0 \int_{(n-\frac{1}{2})\Delta\tau}^{(n+\frac{1}{2})\Delta\tau} d\tau \psi_{\text{R}}(\tau) \quad , \quad n \geq 1 \quad (\text{S8})$$

and, for  $n = 0$ , as

$$F_{\text{R}}^0(0) = N_{\text{R}}^0 \int_0^{\frac{\Delta\tau}{2}} d\tau \psi_{\text{R}}(\tau) \quad (\text{S9})$$

Inserting the Taylor expansion from eq S6 and performing the integrals in eqs S8 and S9, we find

$$F_{\text{R}}^0(n) = N_{\text{R}}^0 \Delta\tau \psi_{\text{R}}(n\Delta\tau) \left[ 1 + \mathcal{O}\left(\frac{\Delta\tau}{\tau_{\text{R}}}\right)^2 \right] \quad , \quad n \geq 1 \quad (\text{S10})$$

and

$$F_{\text{R}}^0(0) = N_{\text{R}}^0 \Delta\tau \psi_{\text{R}}(0) \left[ 1 + \mathcal{O}\left(\frac{\Delta\tau}{\tau_{\text{R}}}\right) \right] \quad (\text{S11})$$

Comparison of eqs S7 and S10 shows that

$$F_{\text{R}}^0(n) = F_{\text{R}}(n) \left[ 1 + \mathcal{O}\left(\frac{\Delta\tau}{\tau_{\text{R}}}\right)^2 \right] \quad , \quad n \geq 1 \quad (\text{S12})$$

Thus, for  $n \geq 1$ ,  $F_{\text{R}}^0(n)$  and  $F_{\text{R}}(n)$  are equal to first order in  $\Delta\tau/\tau_{\text{R}}$ , that is, if we neglect terms of order  $(\Delta\tau/\tau_{\text{R}})^2$ .

To establish the relation between  $N_{\text{R}}$  and  $N_{\text{R}}^0$ , we first note that  $F_{\text{R}}(0)$  is undefined, since we cannot observe less than one frame. It follows immediately from eqs S8 and S9 that

$$\sum_{n=0}^{\infty} F_{\text{R}}^0(n) = N_{\text{R}}^0 \quad (\text{S13})$$



showing that no RTs are lost in this mapping. Thus,

$$N_{\text{R}}^0 = F_{\text{R}}^0(0) + \sum_{n=1}^{\infty} F_{\text{R}}^0(n) \approx F_{\text{R}}^0(0) + \sum_{n=1}^{\infty} F_{\text{R}}(n) = F_{\text{R}}^0(0) + N_{\text{R}} \quad (\text{S14})$$

where, according to eq S12, the second equality holds to first order in  $\Delta\tau/\tau_{\text{R}}$  and the last equality follows from eq S1. With eq S11, we then obtain to first order in  $\Delta\tau/\tau_{\text{R}}$

$$N_{\text{R}} = N_{\text{R}}^0 [1 - \Delta\tau \psi_{\text{R}}(0)] \quad (\text{S15})$$

Since  $\psi_{\text{R}}(0)$  is of order  $1/\tau_{\text{R}}$  (at least for a Poisson process), we see that  $N_{\text{R}}$  and  $N_{\text{R}}^0$  differ to first order in  $\Delta\tau/\tau_{\text{R}}$ , whereas  $F_{\text{R}}(n)$  and  $F_{\text{R}}^0(n)$  for  $n \geq 1$  differ only to second order in  $\Delta\tau/\tau_{\text{R}}$ .

We now consider the binning error in the RT and ST statistics. The trajectory length  $T = N_{\text{F}} \Delta\tau$  can be expressed either in terms of the observed number  $N_{\text{R}}$  of discrete RTs or in terms of the number  $N_{\text{R}}^0$  of continuous RTs:

$$T = N_{\text{F}} \Delta\tau = N_{\text{R}} \tau_{\text{R}} = N_{\text{R}}^0 \tau_{\text{R}}^0 \quad (\text{S16})$$

where  $N_{\text{F}}$  is the number of frames or bins in the trajectory. The frames are labeled  $k = 1, 2, \dots, N_{\text{F}}$ . There are then  $N_{\text{F}} - 1$  intervals  $\Delta\tau$ . But the bins used in connection with  $F_{\text{R}}^0(n)$  include the  $n = 0$  bin, so they also number  $N_{\text{F}}$ .

The mean of the continuous RTs is

$$\tau_{\text{R}}^0 = \Delta\tau \frac{N_{\text{F}}}{N_{\text{R}}^0} \quad (\text{S17})$$

But the mean RT that we compute from the observed RTs is

$$\tau_R = \Delta\tau \frac{N_F}{N_R} \quad (\text{S18})$$

Consequently

$$\tau_R = \tau_R^0 \frac{N_R^0}{N_R} = \frac{\tau_R^0}{1 - \Delta\tau \psi_R(0)} \quad (\text{S19})$$

where eq S15 was used in the last step. This is a systematic error:  $\tau_R$  is slightly longer than the true mean RT  $\tau_R^0$  because RTs shorter than  $\Delta\tau$  may escape detection at a resolution of  $\Delta\tau$ . The mean RT  $\tau_R$  computed from eq S18 is therefore accurate only to zeroth order in  $\Delta\tau/\tau_R$ , that is, it is too long by a relative amount of order  $\Delta\tau/\tau_R$ .

The mean continuous ST can be computed from the identity

$$\tau_S^0 = \frac{\langle \tau^2 \rangle_R^0}{2 \tau_R^0} = \frac{N_R^0}{2 N_F \Delta\tau} \int_0^\infty d\tau \tau^2 \psi_R(\tau) \quad (\text{S20})$$

where eq S17 was used in the second step. By substituting the Taylor expansion of  $\psi_R(\tau)$  from eq S6 into eq S20 and carrying out the integrals, we find to first order in  $\Delta\tau/\tau_R$ , that is, by neglecting terms of order  $(\Delta\tau)^2/(\tau_R \tau_S)$  or higher,

$$\tau_S = \frac{\Delta\tau}{2 N_F} \sum_{n=1}^{\infty} n^2 F_R(n) \quad (\text{S21})$$

Because this expression does not involve  $N_R$ , the mean ST  $\tau_S$  computed from eq S21 is, like  $F_R(n)$  for  $n \geq 1$ , accurate to first order in  $\Delta\tau/\tau_R$ . Discretization introduces a

second-order error in both  $\tau_R$  and  $\tau_S$ . However, the resolution error is of first order in  $\tau_R$  but only of second order in  $\tau_S$ . The resolution error (loss of short RTs) impacts less on  $\tau_S$  than on  $\tau_R$  because short RTs contribute less to  $\tau_S$ .

The continuum residence correlation function (RCF)  $Q_R^0(\tau)$  is defined as

$$Q_R^0(\tau) = \int_{\tau}^{\infty} d\tau' \psi_R(\tau') \quad (\text{S22})$$

Thus,  $Q_R^0(\tau)$  is the fraction of continuous RTs that are longer than  $\tau$ . The discrete RCF  $Q_R^0(n)$  is defined as the fraction of bin-based RTs (of length  $n \Delta\tau$ , including the possibility of  $n = 0$ ) that are longer than  $n \Delta\tau$ . Equation S22 yields to first order

$$Q_R^0(n) = \frac{1}{N_R^0} \sum_{p=n+1}^{\infty} F_R(p) \quad (\text{S23})$$

where we have used eq S12 to replace  $F_R^0(n)$  by  $F_R(n)$ , which is also accurate to first order. But the RCF that we compute is

$$Q_R(n) = \frac{1}{N_R} \sum_{p=n+1}^{\infty} F_R(p) \quad (\text{S24})$$

which has  $N_R$  rather than  $N_R^0$ . Therefore, in view of eq S15,

$$Q_R(n) = \frac{Q_R^0(n)}{1 - \Delta\tau \psi_R(0)} \quad (\text{S25})$$

showing that  $Q_R(n)$  is only accurate to zeroth order in  $\Delta\tau/\tau_R$ .

The continuum survival correlation function (SCF)  $Q_S^0(\tau)$  is related to the RCF via the identity

$$\frac{d Q_S^0(\tau)}{d\tau} = - \frac{Q_R^0(\tau)}{\tau_R^0} \quad (\text{S26})$$

which integrates to

$$Q_S^0(\tau) = \frac{1}{\tau_R^0} \int_{\tau}^{\infty} d\tau' Q_R^0(\tau') \quad (\text{S27})$$

The discrete version of this expression, accurate to first order, is

$$Q_S^0(n) = \frac{N_R^0}{N_F} \sum_{p=n}^{\infty} Q_R^0(p) \quad (\text{S28})$$

where we have used eq S17. If we now substitute  $Q_R^0(p)$  from eq S23, the unknown quantity  $N_R^0$  cancels out and we are left with

$$Q_S^0(n) = \frac{1}{N_F} \sum_{p=n}^{\infty} \sum_{q=p+1}^{\infty} F_R(q) \quad (\text{S29})$$

This expression, which only involves the known quantities  $N_F$  and  $F_R(q)$ , is also the one used to compute the approximate  $Q_S(n)$  from the MD data. We thus see that  $Q_S(n)$  is accurate to first order, whereas  $Q_R(n)$  is only accurate to zeroth order.

## S5. STATISTICAL ERROR

Here we provide more details about the analysis of the statistical uncertainties, or standard deviations,  $\sigma(\tau_X)$  and  $\sigma(Q_X)$  in  $\tau_X$  and  $Q_X$  (with  $X = R$  or  $S$ ), resulting from the finite length  $T$  of the trajectory at a given sampling resolution  $\Delta\tau$ . Here, we assume that  $\Delta\tau$  is so small that the binning error is negligible.

Let  $\tau_X(T)$  be the value of  $\tau_X$  that we compute from a trajectory of length  $T$ . This value generally deviates somewhat from the value  $\tau_X(T \rightarrow \infty)$  that we would compute from an infinitely long trajectory. The standard deviation  $\sigma(\tau_X)$  is a measure of this deviation. If the system is ergodic, the infinite trajectory average is equal to the ensemble average,

$$\langle \tau_X \rangle = \lim_{T \rightarrow \infty} \tau_X(T) \quad (\text{S30})$$

The ensemble is an effectively infinite set of finite trajectories. We can choose to work with either of two ensembles. In the F (frame) ensemble, all trajectories contain the same number  $N_F$  of frames and are therefore of the same length  $T = N_F \Delta\tau$ . In the R (residence time) ensemble, all trajectories contain the same number  $N_R$  of residence time (RT) intervals. The trajectory length is then more conveniently expressed as  $T = N_R \tau_R$ . Because  $\tau_R$  fluctuates somewhat among the different trajectories, it follows that the trajectories in the  $N_R$  ensemble are not of precisely the same length.

Our primary data is a sequence of  $N_R$  RT intervals. The number  $N_R$  is thus a measure of the amount of information at our disposal. If we sample the trajectory more densely by decreasing  $\Delta\tau$ , we increase  $N_F$  but we do not gain more information since  $N_R$  is unchanged. The statistical error should therefore depend on  $N_R$ , but not on  $N_F$ . Specifically, if the RTs are mutually independent, we expect the statistical error to be proportional to  $1/\sqrt{N_R}$ . The statistical error must therefore be evaluated in the  $N_R$  ensemble. In the following, we denote R ensemble averages by  $\langle \cdots \rangle$ , whereas F ensemble averages, in the few cases where they occur, are denoted by  $\langle \cdots \rangle_F$ .

Let  $n$  be the RT in units of  $\Delta\tau$ . The ensemble-averaged mean RT is then

$$\langle n \rangle = \sum_{n=1}^{\infty} n F_R(n) \quad (\text{S31})$$

where  $F_R(n)$  is the equilibrium RT probability distribution. Neither  $F_R(n)$  nor  $\langle n \rangle$  can be obtained from a trajectory of finite length. From the trajectory, we obtain a chronological series of a finite number  $N_R$  of RTs. If the simulated system is in equilibrium, the series  $\{n_\alpha\}_{\alpha=1\dots N_R}$  represents a stationary stochastic process. Stationarity implies, among other things, that the  $n_\alpha$  are identically distributed, but not necessarily independent, random variables. Therefore, the mean  $\langle n_\alpha \rangle \equiv \langle n \rangle$  does not depend on  $\alpha$ . An unbiased estimator of the mean RT  $\langle n \rangle$  is provided by the

trajectory average

$$\bar{n}(N_R) = \frac{1}{N_R} \sum_{\alpha=1}^{N_R} n_{\alpha} \quad (\text{S32})$$

This estimator is unbiased because

$$\langle \bar{n}(N_R) \rangle = \frac{1}{N_R} \sum_{\alpha=1}^{N_R} \langle n_{\alpha} \rangle = \langle n \rangle \quad (\text{S33})$$

Note that  $N_R$  is not ensemble-averaged since it is a constant in the R ensemble. To simplify the notation, we will write  $\bar{n}$  instead of  $\bar{n}(N_R)$  in the following.

The trajectory average  $\bar{n}$  may be regarded as a random variable with a probability density  $\Psi(\bar{n})$ . If  $N_R$  is large, the central limit theorem leads us to expect that  $\Psi(\bar{n})$  is approximately Gaussian and hence fully characterized by the mean  $\langle \bar{n} \rangle = \langle n \rangle$  and the variance

$$\sigma^2(\bar{n}) \equiv \langle [\bar{n} - \langle \bar{n} \rangle]^2 \rangle = \langle \bar{n}^2 \rangle - \langle \bar{n} \rangle^2 \quad (\text{S34})$$

The desired measure of the statistical error in the estimator  $\bar{n}$  of the mean RT is the standard deviation  $\sigma(\bar{n})$ . However, eq S34 involves ensemble averages, which we cannot obtain from a finite trajectory. We can only obtain an estimator  $s^2(\bar{n})$  of the variance  $\sigma^2(\bar{n})$ . To be reliable, this should be an unbiased estimator.

Substituting  $\bar{n}$  from eq S32 into eq S34, we obtain

$$\sigma^2(\bar{n}) = \frac{1}{N_R^2} \sum_{\alpha=1}^{N_R} \sum_{\beta=1}^{N_R} C_{\alpha\beta} \quad (\text{S35})$$

with the covariance  $C_{\alpha\beta}$  of the random variables  $n_\alpha$  and  $n_\beta$  defined as

$$C_{\alpha\beta} \equiv \langle n_\alpha n_\beta \rangle - \langle n_\alpha \rangle \langle n_\beta \rangle = \langle n_\alpha n_\beta \rangle - \langle n \rangle^2 \quad (\text{S36})$$

If the random variables  $n_\alpha$  and  $n_\beta$  are uncorrelated, we have  $\langle n_\alpha n_\beta \rangle = \langle n_\alpha \rangle \langle n_\beta \rangle = \langle n \rangle^2$  and, therefore,  $C_{\alpha\beta} = 0$ . Moreover,  $C_{\alpha\alpha} = \sigma_n^2$  is the variance of the random variable  $n_\alpha$ ,

$$\sigma_n^2 \equiv \langle n_\alpha^2 \rangle - \langle n_\alpha \rangle^2 = \langle n^2 \rangle - \langle n \rangle^2 \quad (\text{S37})$$

In general,  $|C_{\alpha\beta}| \leq \sigma_n^2$ .

It is convenient to split  $\sigma^2(\bar{n})$  in two parts, corresponding to the diagonal and off-diagonal terms of the double sum in eq S35. Thus,

$$\sigma^2(\bar{n}) = \sigma_0^2(\bar{n}) + \Delta_c(\bar{n}) \quad (\text{S38})$$

$$\sigma_0^2(\bar{n}) = \frac{\sigma_n^2}{N_R} \quad (\text{S39})$$

$$\Delta_c(\bar{n}) = \frac{1}{N_R^2} \sum_{\alpha=1}^{N_R} \sum_{\beta=1}^{N_R}{}' C_{\alpha\beta} \quad (\text{S40})$$

where the prime on the second sum signifies omission of diagonal ( $\beta = \alpha$ ) terms.

Whereas  $\sigma^2(\bar{n})$  and  $\sigma_0^2(\bar{n})$  are necessarily non-negative, the quantity  $\Delta_c$  may be negative.

If all the random variables  $\{n_\alpha\}_{\alpha=1\dots N_R}$  are mutually uncorrelated, it follows from the foregoing that  $\Delta_c(\bar{n}) = 0$  and that  $\sigma^2(\bar{n}) = \sigma_0^2(\bar{n})$  with

$$\sigma_0^2(\bar{n}) = \frac{\sigma_n^2}{N_R} = \frac{\langle n^2 \rangle - \langle n \rangle^2}{N_R} \quad (\text{S41})$$



This is a well-known result: the standard error of the mean of  $N$  independent measurements of a random variable equals the standard deviation of the random variable itself divided by the square root of the number of measurements.

To obtain an unbiased estimator  $s_0^2(\bar{n})$  of  $\sigma_0^2(\bar{n})$ , we need an unbiased estimator  $s_n^2$  of the variance  $\sigma_n^2$ . This is

$$s_n^2 = \frac{1}{(N_R - 1)} \sum_{\alpha=1}^{N_R} (n_\alpha^2 - \bar{n}^2) \quad (\text{S42})$$

with the denominator  $N_R - 1$  rather than  $N_R$ . In the absence of correlations, the unbiased estimator  $s_0^2(\bar{n})$  of the variance  $\sigma_0^2(\bar{n})$  of the unbiased estimator  $\bar{n}$  of the mean RT  $\langle n \rangle$  can, according to eqs S32, S41 and S42, be computed from

$$s_0^2(\bar{n}) = \frac{1}{(N_R - 1)} \left[ \left( \frac{1}{N_R} \sum_{\alpha=1}^{N_R} n_\alpha^2 \right) - \left( \frac{1}{N_R} \sum_{\alpha=1}^{N_R} n_\alpha \right)^2 \right] \quad (\text{S43})$$

The desired statistical uncertainty in  $\tau_R/\Delta\tau$  is the square root of this.

We now return to the general case with correlations. If the set  $\{n_\alpha\}_{\alpha=1\dots N_R}$  of random variables represents a stochastic process with the index  $\alpha$  increasing monotonically with time, then the covariance  $C_{\alpha\beta}$  in eq S36 can be regarded as a serial correlation function. If the process is stationary, the correlation function can only depend on the difference  $\gamma = \beta - \alpha$ . The double sum in eq S40 can therefore be reduced to a single sum as

$$\Delta_c(\bar{n}) = \frac{2}{N_R} \sum_{\gamma=1}^{N_R-1} \left( 1 - \frac{\gamma}{N_R} \right) C_\gamma \quad (\text{S44})$$

$$C_\gamma \equiv \langle n_\alpha n_{\alpha+\gamma} \rangle - \langle n_\alpha \rangle \langle n_{\alpha+\gamma} \rangle = \langle n_\alpha n_{\alpha+\gamma} \rangle - \langle n \rangle^2 \quad (\text{S45})$$

where  $\langle n_\alpha n_{\alpha+\gamma} \rangle$  is independent of  $\alpha$ . In our case, the ‘time’ point  $\alpha$  is simply the serial number of a RT in the chronological RT list. Hence, we refer to  $C_\gamma$  as a serial correlation function rather than a time correlation function. The serial correlation function  $C_\gamma$  measures the correlation of two RTs that are separated by  $\gamma - 1$  intervening RTs.

Without *a priori* knowledge about the serial correlation function  $C_\gamma$ , eq S44 is of limited utility. Probably the most robust and computationally efficient approach to this problem is the block renormalization (BR) method proposed by Flyvbjerg and Petersen.<sup>1</sup> In this method, we select the longest sequence of  $N = 2^k$  ( $k$  is an integer) RTs contained in the original RT series. For this pruned series, we compute the  $\sigma(\bar{n})$  estimator from eq S43 (which is valid in the absence of correlations)

$$s = \left\{ \frac{1}{(N-1)} \left[ \left( \frac{1}{N} \sum_{\alpha=1}^N n_\alpha^2 \right) - \left( \frac{1}{N} \sum_{\alpha=1}^N n_\alpha \right)^2 \right] \right\}^{1/2} \quad (\text{S46})$$

We then perform the first blocking transformation by averaging pairwise values

$$n'_\alpha = \frac{n_{2\alpha-1} + n_{2\alpha}}{2} \quad (\text{S47})$$

so that  $N' = N/2$ . We then compute a new value for the estimator,

$$s' = \left\{ \frac{1}{(N'-1)} \left[ \left( \frac{1}{N'} \sum_{\alpha=1}^{N'} n'^2_\alpha \right) - \left( \frac{1}{N'} \sum_{\alpha=1}^{N'} n'_\alpha \right)^2 \right] \right\}^{1/2} \quad (\text{S48})$$

This process is repeated until  $N' = 2$ , yielding a series of  $k$  values for the estimator. The relative standard deviation in these estimators is  $[2(N' - 1)]^{-1/2}$ .<sup>1</sup> Whereas the mean  $\bar{n}$  is invariant under the blocking transformation, the standard deviation estimator  $s$  increases until a so-called fixed point is reached and then remains constant upon further transformations. Of course, such a plateau will only be observed if the correlation ‘time’ on which  $C_\gamma$  decays is much smaller than  $N_R$ .

Figure S2 shows estimates of  $\sigma(\tau_R)$  after successive blocking transformations for site W113 in all states. Here,  $N_R = 37,499$ , so that  $k = 15$  ( $2^{15} = 32,768$ ). Equation S43 yields  $s_0(\tau_R) \approx 2$  ns, but the BR estimate increases to  $\sim 10$  ns without reaching a well-defined plateau. The source of the problem is evident from the RT vector in Fig. 2 of the main text: small and large RTs tend to cluster in alternating segments with a correlation ‘time’ of the same order of magnitude as  $N_R$ . In Fig. S3, the BR method is applied to the RTs of site W113 in state M1 after omission of short RTs. Here the problem is the small number of RTs,  $N_R = 60$ , only half of which can be used in the BR method ( $2^5 = 32$ ). Because of this reduction in the number of RTs, the initial error in Fig. S3 is larger than the one computed from eq S43 with all 60 RTs. Although the uncertainty in  $s(\tau_R)$  is large, the BR method indicates that correlations are unimportant, as is also suggested by the RT vector in Fig. 6 of the main text. The uncorrelated error  $s_0(\tau_R)$  should then be a good estimate.

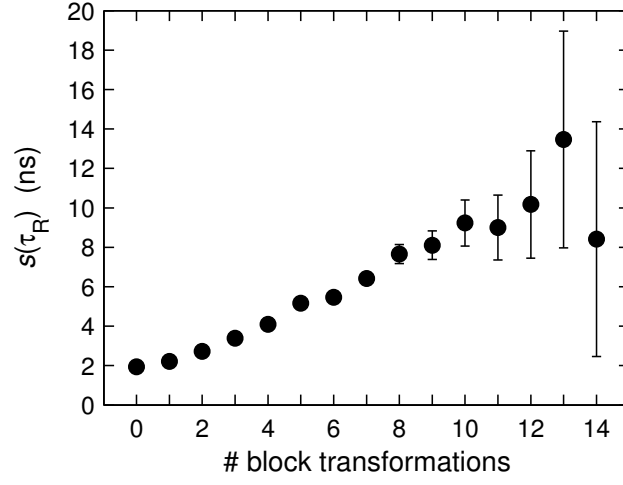


Figure S2: BR estimation of statistical error in  $\tau_R$  for site W113 in all states.

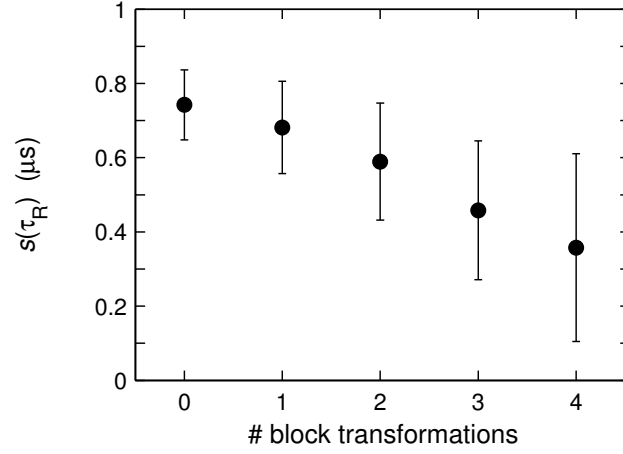


Figure S3: BR estimation of statistical error in  $\tau_R$  for site W113 in M1 state with exclusion of short ( $< 50$  ns) RTs.

We now turn to the mean ST. Let  $m_k$  be the number of remaining frames in the current RT, of length  $n_\alpha$ , beginning at, and including, an arbitrary frame  $k$ . This is the time to the next exchange event from an arbitrary starting time. From a finite-length trajectory, we obtain a chronological series of  $N_F$  STs. If the simulated system is in equilibrium, the series  $\{m_k\}_{k=1\dots N_F}$  represents a stationary stochastic process and the trajectory average

$$\overline{m} = \frac{1}{N_F} \sum_{k=1}^{N_F} m_k \quad (\text{S49})$$

is an unbiased estimator of the mean ST  $\langle m \rangle_F$  in the F ensemble.

Since the RT series  $\{n_\alpha\}_{\alpha=1\dots N_R}$  contains all information about the underlying stochastic process, it must be possible to express  $\overline{m}$  in terms of the RTs  $n_\alpha$ . Indeed, it is readily shown that (see eq 7 of the main text)

$$\overline{m} = \frac{\overline{n^2}}{2\overline{n}} = \frac{1}{2} \frac{\sum_{\alpha=1}^{N_R} n_\alpha^2}{\sum_{\alpha=1}^{N_R} n_\alpha} \quad (\text{S50})$$

We seek the variance

$$\sigma^2(\overline{m}) \equiv \langle [\overline{m} - \langle \overline{m} \rangle]^2 \rangle = \langle \overline{m}^2 \rangle - \langle \overline{m} \rangle^2 \quad (\text{S51})$$

We can choose to evaluate the ensemble averages either in the F ensemble or in the R ensemble. In the former case, we can substitute  $\overline{m}$  from eq S49 and proceed as for  $\tau_R$ . But we want an expression that involves the RTs  $n_\alpha$ , since the RT series

constitutes our data. We must therefore work with the R ensemble, where  $N_R$  is constant. Our starting point must therefore be eq S50, which expresses the random variable  $\overline{m}$  as a nonlinear function of two other (correlated) random variables  $\overline{n}$  and  $\overline{n^2}$ . The conventional procedure for computing the variance of a function of random variables is to Taylor expand the function. Expanding eq S50 to first order, we obtain

$$\overline{m} = \langle m \rangle - \frac{\langle n^2 \rangle}{2 \langle n \rangle^2} [\overline{n} - \langle n \rangle] + \frac{1}{2 \langle n \rangle} [\overline{n^2} - \langle n^2 \rangle] + \dots \quad (\text{S52})$$

where we have noted that  $\langle \overline{n} \rangle = \langle n \rangle$  and  $\langle \overline{n^2} \rangle = \langle n^2 \rangle$  in the R ensemble. This expansion shows that  $\langle \overline{m} \rangle = \langle m \rangle$  to first order in the R ensemble, whereas this is an exact result in the F ensemble.

Forming the difference  $\overline{m} - \langle m \rangle$  from eq S52, squaring it, and taking the ensemble average, we obtain the variance

$$\begin{aligned} \sigma^2(\overline{m}) = & \left( \frac{\langle n^2 \rangle}{2 \langle n \rangle^2} \right)^2 \left\langle [\overline{n} - \langle n \rangle]^2 \right\rangle + \left( \frac{1}{2 \langle n \rangle} \right)^2 \left\langle [\overline{n^2} - \langle n^2 \rangle]^2 \right\rangle \\ & - \frac{\langle n^2 \rangle}{2 \langle n^3 \rangle} \left\langle [\overline{n} - \langle n \rangle] [\overline{n^2} - \langle n^2 \rangle] \right\rangle \end{aligned} \quad (\text{S53})$$

Since we have neglected terms of higher than second order, this result is accurate only if the relative variances in  $\overline{n}$  and  $\overline{n^2}$ , as well as their relative covariance, are small. In other words, the Taylor expansion method only predicts  $\sigma(\overline{m})$  accurately if  $\sigma(\overline{m})/\overline{m}$  is small. Furthermore, the estimator obtained by replacing all ensemble averages in eq S53 by trajectory averages is not unbiased.

Rather than pursuing the general case, we shall assume that the RTs  $n_\alpha$  are mutually uncorrelated. This seems to be approximately true if we omit the short RTs. In any case,  $\sigma_0^2(\overline{m})$  provides a useful lower bound on  $\sigma^2(\overline{m})$ . For uncorrelated RTs, we have

$$\begin{aligned}\langle [\overline{n} - \langle n \rangle]^2 \rangle &= \frac{\langle n^2 \rangle - \langle n \rangle^2}{N_R} \\ \langle [\overline{n^2} - \langle n^2 \rangle]^2 \rangle &= \frac{\langle n^4 \rangle - \langle n^2 \rangle^2}{N_R} \\ \langle [\overline{n} - \langle n \rangle] [\overline{n^2} - \langle n^2 \rangle] \rangle &= \frac{\langle n^3 \rangle - \langle n \rangle \langle n^2 \rangle}{N_R}\end{aligned}\tag{S54}$$

which is inserted into eq S53 to yield

$$\sigma_0^2(\overline{m}) = \frac{1}{N_R} \frac{\langle n \rangle^2 \langle n^4 \rangle + \langle n^2 \rangle^3 - 2 \langle n \rangle \langle n^2 \rangle \langle n^3 \rangle}{4 \langle n \rangle^4}\tag{S55}$$

As expected, the error  $\sigma_0(\overline{m})$  decreases as the inverse square root of  $N_R$ , just as for  $\sigma_0(\overline{n})$ . However, the estimator  $s_0^2(\overline{m})$  obtained by replacing the ensemble averages in eq S55 by trajectory averages (and perhaps  $N_R$  by  $N_R - 1$ ) is not unbiased. Moreover, we cannot use the BR method to obtain  $\sigma^2(\overline{m})$  (including the effect of correlations), because the trajectory average  $\overline{m}$  in eq S50 is not invariant under blocking transformation. Nevertheless, we can hope that eq S55 yields a reasonable estimate if the relative error is small and if the correlations are insignificant.

If the  $n_\alpha$  are not only mutually independent but also exponentially distributed,

they constitute a Poisson process and

$$\langle n^k \rangle = k! \langle n \rangle^k \quad (\text{S56})$$

In this special case, eq S55 reduces to

$$\sigma_0^2(\overline{m}) = \frac{2 \langle n \rangle^2}{N_R} \quad (\text{S57})$$

whereas the corresponding result in eq S41 for the variance of the mean RT becomes

$$\sigma_0^2(\overline{n}) = \frac{\langle n \rangle^2}{N_R} \quad (\text{S58})$$

Whereas  $\tau_S = \tau_R$  (or  $\overline{m} = \overline{n}$ ) for a Poisson process, the standard deviation is thus a factor  $\sqrt{2}$  larger for  $\tau_S$ .

The trajectory-based estimator of the residence correlation function (RCF)  $Q_R(n)$  is denoted by  $\overline{h}(n)$ . This may be expressed as a trajectory average

$$\overline{h}(n) = \frac{1}{N_R} \sum_{\alpha=1}^{N_R} h_{\alpha}(n) \quad (\text{S59})$$

of the indicator function  $h_{\alpha}(n)$ , which equals 1 if  $n_{\alpha} > n$  and 0 otherwise. The ensemble-averaged RCF is

$$\langle h(n) \rangle = \sum_{p=n+1}^{\infty} F_R(p) \quad (\text{S60})$$

where  $F_R(p)$  is the equilibrium RT probability distribution. The trajectory average  $\overline{h}(n)$  is an unbiased estimator of the ensemble average  $\langle h(n) \rangle$ , because

$$\langle \overline{h}(n) \rangle = \frac{1}{N_R} \sum_{\alpha=1}^{N_R} \langle h_{\alpha}(n) \rangle = \langle h(n) \rangle \quad (\text{S61})$$



where we have made use of the stationarity of the stochastic process  $\{n_\alpha\}_{\alpha=1\dots N_R}$ .

The error analysis for  $\tau_R$  carries over directly to the RCF if we replace  $n_\alpha$  by  $h_\alpha(n)$ . To simplify the notation, we shall omit the argument  $n$  in most of what follows. Because  $h_\alpha$  can only take the values 0 or 1, the expressions can be further simplified by using the identity  $h_\alpha^2 = h_\alpha$ . The variance of  $h_\alpha$  is thus given by

$$\sigma_h^2 \equiv \langle h_\alpha^2 \rangle - \langle h_\alpha \rangle^2 = \langle h \rangle [1 - \langle h \rangle] \quad (\text{S62})$$

If the RTs are mutually uncorrelated, the variance of  $\bar{h}$  is, in analogy with eq S41,

$$\sigma_0^2(\bar{h}) = \frac{\sigma_h^2}{N_R} = \frac{\langle h \rangle [1 - \langle h \rangle]}{N_R} \quad (\text{S63})$$

and the unbiased estimator (see eq S43) for the standard deviation is

$$s_0(\bar{h}) = \left[ \frac{\bar{h}(1 - \bar{h})}{N_R - 1} \right]^{1/2} \quad (\text{S64})$$

Since  $\bar{h}(0) \equiv 1$ , eq S64 predicts a vanishing error for  $n = 0$ , as expected.

The trajectory-based estimator of the survival correlation function (SCF)  $Q_S(n)$  is denoted by  $\bar{g}(n)$ . This may be expressed as a trajectory average

$$\bar{g}(n) = \frac{1}{N_F} \sum_{k=1}^{N_F} g_k(n) \quad (\text{S65})$$

of the indicator function  $g_k(n)$ , which equals 1 if  $m_k > n$  and 0 otherwise.

The ensemble-averaged SCF is

$$\langle g(n) \rangle = \sum_{p=n+1}^{\infty} F_S(p) \quad (\text{S66})$$

where  $F_S(p)$  is the equilibrium ST probability distribution. The trajectory average  $\bar{g}(n)$  is an unbiased estimator of the ensemble average  $\langle g(n) \rangle_F$  in the F ensemble, because the stochastic process is stationary so that

$$\langle \bar{g}(n) \rangle_F = \frac{1}{N_F} \sum_{k=1}^{N_F} \langle g_k(n) \rangle_F = \langle g(n) \rangle_F \quad (\text{S67})$$

As for  $\tau_S$ , we want to express the trajectory average in eq S65 in terms of the random variables  $n_\alpha$ , so we can take ensemble averages in the R ensemble. We therefore write

$$\bar{g}(n) = \sum_{\alpha=1}^{N_R} p_\alpha g_\alpha(n) \quad (\text{S68})$$

where  $g_\alpha(n)$  is the mean of the indicator function  $g_k(n)$  for the RT interval  $\alpha$  and  $p_\alpha = n_\alpha/N_F = n_\alpha/(N_R \bar{n})$  is the probability that the randomly chosen frame  $k$  lies in this interval. Furthermore,

$$g_\alpha(n) = \frac{1}{n_\alpha} \sum_{i=1}^{n_\alpha} g_i(n) = \frac{1}{n_\alpha} (n_\alpha - n) h_\alpha(n) \quad (\text{S69})$$

where the indicator function  $h_\alpha(n)$ , introduced in connection with the RCF, acts as a step function that ensures that  $g_\alpha(n) = 0$  for  $n > n_\alpha$ . Combining eqs S68 and S69 and the expression for  $p_\alpha$ , we obtain

$$\bar{g}(n) = \frac{1}{N_R \bar{n}} \sum_{\alpha=1}^{N_R} (n_\alpha - n) h_\alpha(n) \quad (\text{S70})$$

Note that, in eqs S69 and S70,  $n$  is not a random variable, but merely an integer.

We seek the variance

$$\sigma^2(\bar{g}) \equiv \langle [\bar{g} - \langle \bar{g} \rangle]^2 \rangle = \langle \bar{g}^2 \rangle - \langle \bar{g} \rangle^2 \quad (\text{S71})$$

Now eq S70 expresses the random variable  $\bar{g}$  as a nonlinear function of  $2N_R + 1$  other (correlated) random variables, namely  $n_\alpha$ ,  $h_\alpha$  and  $\bar{n}$ . Taylor expanding  $\bar{g}(n)$  in eq S70 to first order in the fluctuations of these variables around their ensemble averages, we can proceed along the same lines as in the error analysis of  $\tau_S$ . However, this leads to a rather complicated result, so we shall pursue a different approach that yields a simpler, albeit approximate, result that is likely to be a slight overestimate of the true (uncorrelated) error in  $Q_S(\tau)$ .

The starting point is eq S27, which we now write as

$$Q_S(\tau) = 1 - \frac{1}{\tau_R} \int_0^\tau d\tau' Q_R(\tau') \quad (\text{S72})$$

The discrete version of this expression is

$$Q_S(n) = 1 - \frac{\Delta\tau}{\tau_R} \sum_{p=0}^{n-1} Q_R(p) = 1 - \frac{\Delta\tau}{\tau_R} \left[ \sum_{p=0}^n Q_R(p) - Q_R(n) \right] \quad (\text{S73})$$

or, in terms of dimensionless trajectory averages,

$$\bar{g}(n) = 1 - \frac{1}{\bar{n}} \left[ \sum_{p=0}^n \bar{h}(p) - \bar{h}(n) \right] \quad (\text{S74})$$

In view of eq S59, this can be expressed as

$$\bar{g}(n) = 1 - \frac{1}{N_R} \sum_{\alpha=1}^{N_R} G_\alpha(n) \quad (\text{S75})$$

where we have defined

$$G_\alpha(n) \equiv \frac{1}{\bar{n}} \left[ \sum_{p=0}^n h_\alpha(p) - h_\alpha(n) \right] \quad (\text{S76})$$

We now introduce the approximation, which consists in neglecting the fluctuations in  $\bar{n}$  so eq S76 becomes a linear combination of random variables,

$$G_\alpha(n) \approx \frac{1}{\langle n \rangle} \left[ \sum_{p=0}^n h_\alpha(p) - h_\alpha(n) \right] \quad (\text{S77})$$

As before, we shall only consider the uncorrelated case, where  $\langle n_\alpha n_\beta \rangle = \langle n_\alpha \rangle \langle n_\beta \rangle$  and therefore  $\langle h_\alpha(n) h_\beta(p) \rangle = \langle h_\alpha(n) \rangle \langle h_\beta(p) \rangle$  and  $\langle G_\alpha(n) G_\beta(p) \rangle = \langle G_\alpha(n) \rangle \langle G_\beta(p) \rangle$ .

It then follows immediately that

$$\sigma_0^2(\bar{g}) = \frac{\sigma_G^2}{N_R} \quad (\text{S78})$$

where  $\sigma_G^2$  is the variance of  $G_\alpha(n)$ ,

$$\sigma_G^2 = \langle [G_\alpha(n)]^2 \rangle - \langle G_\alpha(n) \rangle^2 \quad (\text{S79})$$

Ensemble averaging eq S77, we obtain

$$\langle G_\alpha(n) \rangle = \frac{1}{\langle n \rangle} \left[ \sum_{p=0}^n \langle h(p) \rangle - \langle h(n) \rangle \right] \quad (\text{S80})$$

Similarly,

$$\begin{aligned} \langle [G_\alpha(n)]^2 \rangle &= \frac{1}{\langle n \rangle^2} \left[ \sum_{p=0}^n \sum_{q=0}^n \langle h_\alpha(p) h_\alpha(q) \rangle - 2 \sum_{p=0}^n \langle h_\alpha(n) h_\alpha(p) \rangle + \langle h_\alpha^2(n) \rangle \right] \\ &= \frac{1}{\langle n \rangle^2} \left[ \sum_{p=0}^n (2p+1) \langle h(p) \rangle - (2n+1) \langle h(n) \rangle \right] \end{aligned} \quad (\text{S81})$$

where we have made use of the identities

$$h_{\alpha}^2(n) = h_{\alpha}(n) \quad (\text{S82})$$

$$h_{\alpha}(n) h_{\alpha}(p) = h_{\alpha}(n) \quad \text{if } p \leq n \quad (\text{S83})$$

Combining eqs S78 – S81, we obtain

$$\sigma_0^2(\bar{g}) = \frac{1}{\langle n \rangle^2 N_{\text{R}}} \left\{ \sum_{p=0}^n (2p+1) \langle h(p) \rangle - (2n+1) \langle h(n) \rangle - \left[ \sum_{p=0}^n \langle h(p) \rangle - \langle h(n) \rangle \right]^2 \right\} \quad (\text{S84})$$

The estimator for the standard deviation of  $\bar{g}(n)$  is thus

$$s_0(\bar{g}) = \frac{1}{\bar{n} N_{\text{R}}^{1/2}} \left\{ \sum_{p=0}^n (2p+1) \bar{h}(p) - (2n+1) \bar{h}(n) - \left[ \sum_{p=0}^n \bar{h}(p) - \bar{h}(n) \right]^2 \right\}^{1/2} \quad (\text{S85})$$

As for  $\bar{h}(n)$ , eq S85 shows that the error in  $\bar{g}(n)$  has the expected property of vanishing for  $n = 0$ .

## References

<sup>1</sup> Flyvbjerg, H.; Petersen, H. G. *J. Chem. Phys.* **1989**, *91*, 461 – 466.